# FOURIERROFORMER: LEARNED FOURIER ATTENTION FOR VISION TRANSFORMERS

## Anonymous authors

Paper under double-blind review

## **ABSTRACT**

Vision Transformers (ViTs) excel at long-range reasoning but lack principled mechanisms for modeling spatial frequencies and controlling how attention decays with distance. We propose FourierRoFormer, a frequency-aware Transformer that augments rotary positional embeddings with learnable Fourier components. This enables explicit modeling of multi-scale visual patterns and adaptive distance-dependent modulation of attention. Our analysis shows that FourierRo-Former produces attention hierarchies aligned with object boundaries (correlation r = 0.85) and distinct specialization across attention heads. On ImageNet-1K, FourierRoFormer achieves 84.1% top-1 accuracy (+1.8pp over RoFormer) while using 25% fewer parameters than competitive spectral methods. It also improves transfer to dense prediction tasks, yielding +2.6 mAP on COCO detection and +2.2 mAP on instance segmentation. Ablation studies highlight the complementary roles of frequency modulation (+4.43pp) and adaptive damping (+2.09pp). Despite its expressiveness, the method introduces only 0.04% additional parameters and  $\sim 3\%$  computational overhead, confirmed by complexity and FLOPs analysis.

## 1 Introduction

Transformer architectures have become the dominant paradigm across vision, language, and multimodal learning (Vaswani et al., 2017; Dosovitskiy et al., 2020; Brown et al., 2020). In computer vision, Vision Transformers (ViTs) (Dosovitskiy et al., 2020) have achieved consistent improvements in recognition tasks by treating images as sequences of patches and applying self-attention to capture global dependencies.

However, standard attention mechanisms face key limitations when processing structured visual data: (1) they lack inductive bias about spatial relationships, (2) they are frequency-blind to the multi-scale nature of visual patterns, and (3) they provide limited control over how attention decays across token distances (Park & Kim, 2022; Raghu et al., 2021; Rao et al., 2021; Press et al., 2021). Recent approaches such as relative positional encodings (Shaw et al., 2018), rotary embeddings (Su et al., 2024), and windowed attention (Liu et al., 2021) improve spatial awareness but still fall short of explicitly modeling frequency relationships.

We address these challenges by drawing on principles from signal processing and propose **FourierRoFormer**. Our method integrates learnable Fourier components into the transformer attention mechanism, enabling frequency-aware modulation of attention scores as a function of token distance. Unlike prior rotary or Fourier-based models, FourierRoFormer adaptively learns which frequency bands are most relevant for visual understanding. Figure 1 illustrates how Fourier modulation reshapes attention to emphasize multi-scale structures. This perspective provides a principled way to control information propagation across scales, bridging the gap between spectral theory and transformer design.

By incorporating a learnable mixture of sinusoidal components with frequencies, amplitudes, and phases, FourierRoFormer adaptively modulates attention based on token distances (Section 3). Our unified framework combines Fourier modulation with rotary positional embeddings and optional exponential damping. Theoretical analysis explains how these components influence attention gradients and feature propagation (Appendix A). Extensive experiments demonstrate that FourierRoFormer consistently outperforms ViT, DeiT, and RoFormer baselines, while ablations high-

light the complementary effects of frequency modulation and damping, providing insights into how frequency-aware attention improves multiscale feature capture (Section 4). These contributions establish FourierRoFormer as a principled framework for frequency-aware Transformers.

## 2 Related Work

The *Vision Transformer* (ViT) (Dosovitskiy et al., 2020) was the first to show that the transformer architecture—originally designed for language—can excel at image classification by cutting images into fixed-size patches and treating each as a token for self-attention. Although ViT achieves strong accuracy on large datasets, it requires much more training data than traditional convolutional networks. Follow-up work like DeiT (Touvron et al., 2021) addressed this data-hunger with distillation and augmentation, while Swin (Liu et al., 2021) and PVT (Wang et al., 2021) introduced hierarchical, multi-scale designs (shifted windows in Swin; a pyramid with spatial-reduction attention in PVT). In parallel, spectral token-mixing approaches leverage fixed transforms in the frequency domain—Fourier, wavelet, or scattering—either to replace or to augment attention (e.g., GFNet, Wave-ViT, SpectFormer, SVT) (Rao et al., 2021; Yao et al., 2022; Patro et al., 2025a; Patro & Agneeswaran, 2023). While standard dot-product attention is not explicitly frequency-aware, spectral components inject frequency-selective inductive bias that is complementary to hierarchical and locality biases. In this work, we introduce *FourierRoFormer*, which aims to address this frequency-blindness by embedding frequency-aware modulation directly into the attention scores.

Beyond the challenge of frequency awareness, transformers face another fundamental limitation: self-attention is permutation-invariant, so transformers need an additional signal to recover token order (Vaswani et al., 2017). RoPE (Su et al., 2024) rotates query and key vectors, so their inner product encodes relative distance, but still treats all frequencies uniformly with no control over attention decay. FourierRoFormer extends RoPE by learning sinusoid mixtures whose parameters are data-optimized, providing interpretable frequency-selective attention decay.

Several studies speed up attention by approximating its  $\mathcal{O}(n^2)$  complexity. Performer (Choromanski et al., 2020) and Linformer (Wang et al., 2020) use low-rank projections; EfficientFormer (Li et al., 2022) and MobileViT (Mehta & Rastegari, 2021) redesign the backbone for mobile deployment. These methods mainly target runtime and memory, leaving the *frequency content* of attention untouched. In contrast, FourierRoFormer focuses on richer signal modeling while retaining a compute profile comparable to standard RoPE attention.

Complementing these efficiency-focused approaches, there is growing interest in incorporating frequency analysis principles into neural networks. Frequency analysis has deep roots in signal processing and is increasingly common in modern networks.

## 3 METHODOLOGY

In this section, we present the FourierRoFormer architecture, detailing how Fourier components are integrated into the attention mechanism and describing the overall model design (Figure 1). Detailed mathematical analyses, proofs, and additional properties are provided in the appendices.

To establish the foundation for our approach, we first briefly review standard transformer attention. Standard transformer self-attention (Vaswani et al., 2017) computes attention scores between query  $\mathbf{Q} \in \mathbb{R}^{n \times d}$  and key  $\mathbf{K} \in \mathbb{R}^{n \times d}$  matrices as  $\mathbf{A} = \operatorname{softmax} \left( \mathbf{Q} \mathbf{K}^{\top} / \sqrt{d} \right)$  These weights compute a weighted sum of value vectors  $\mathbf{V} \in \mathbb{R}^{n \times d}$  using Attention $(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \mathbf{A} \mathbf{V}$ . This formulation treats all token interactions uniformly regardless of their spatial relationships, a key limitation for visual data with multiscale patterns. While RoPE (Su et al., 2024) partially addresses this by encoding relative positions through rotations:  $\langle \mathbf{q}_m^{\text{RoPE}}, \mathbf{k}_n^{\text{RoPE}} \rangle = \langle \mathbf{R}_{\theta,m} \mathbf{q}_m, \mathbf{R}_{\theta,n} \mathbf{k}_n \rangle$ , it still lacks explicit frequency awareness (further analysis is in Appendix D).

Building upon RoPE's relative positioning capabilities, we introduce FourierRoFormer attention, which enhances RoPE with learned Fourier modulation and optional damping, as illustrated in Figure 1.

**Fourier Modulation.** The Fourier modulation function  $\mathcal{M}(d)$  is defined as a weighted sum of cosine functions with learnable frequencies, amplitudes, and phases (see Figure 1, bottom panel):

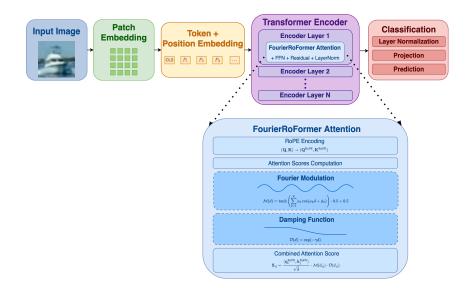


Figure 1: FourierRoFormer architecture for Vision Transformers. The figure illustrates the complete pipeline from input image through patch and position embeddings, transformer encoder, to classification head. It also details the FourierRoFormer attention mechanism, featuring RoPE mappings, attention score computation, Fourier modulation via learnable components  $\mathcal{M}(d)$ , and exponential damping  $\mathcal{D}(d)$  to control distance-based decay.

$$\mathcal{M}(d) = \frac{1}{2} \left( \tanh \left( \sum_{k=1}^{K} a_k \cos(\omega_k d + \phi_k) \right) + 1 \right)$$
 (1)

where K is the number of Fourier components,  $a_k$  are the amplitudes,  $\omega_k$  are the frequencies, and  $\phi_k$  are the phase shifts. The hyperbolic tangent and scaling ensure that the modulation values lie in the range (0,1), allowing the model to attenuate attention scores in a continuous manner. Unlike standard attention, which treats all token pairs identically, our formulation learns which frequency patterns are most relevant for visual tasks. High-frequency components capture fine-grained details while low-frequency components maintain global context—precisely the multi-scale capability standard attention lacks. A detailed analysis of the properties of this modulation function, including its approximation capabilities and interpretability, is presented in Appendix A.

**Proposition 1** (Interpretability of Fourier Components). For each basis element in modulation function  $\mathcal{M}(d)$ , amplitude  $a_k$  dictates how strongly the k-th cosine term contributes—the larger  $|a_k|$ , the greater its influence. Frequency  $\omega_k$  sets the spatial oscillation rate; higher values produce finergrained overall attention variation as token distance d changes. Finally, phase shift  $\phi_k$  translates the component horizontally along the distance axis, relocating attention peaks and troughs while leaving frequency intact.

This formulation enables the model to learn periodic patterns for modulating attention by token distance. By mixing sinusoidal components, it captures multi-scale relationships, selectively emphasizing or suppressing connections based on data characteristics. Unlike standard attention mechanisms, our approach learns which frequency patterns are most relevant for visual understanding. The theoretical foundation for this modulation function is established by the following key theorem, which demonstrates its approximation capabilities and interpretability properties:

**Theorem 1** (Properties of Fourier Modulation Function). Let  $\mathcal{M}: \mathbb{R} \to (0,1)$  be the Fourier modulation function defined in equation I, where  $a_k \in \mathbb{R}$  are learnable amplitudes,  $\omega_k > 0$  are learnable frequencies, and  $\phi_k \in [0,2\pi)$  are learnable phase shifts for  $k=1,\ldots,K$ . Then  $\mathcal{M}(d)$  is a smooth function with  $\mathcal{M}(d) \in (0,1)$  for all  $d \in \mathbb{R}$ . For any continuous function  $f:[0,L] \to (0,1)$  and any  $\varepsilon > 0$ , there exists an integer K and parameters  $\{a_k, \omega_k, \phi_k\}_{k=1}^K$  such that

$$\sup_{d \in [0,L]} |\mathcal{M}(d) - f(d)| < \varepsilon.$$

If the set of frequencies  $\{\omega_k\}_{k=1}^K$  consists of rational multiples of each other, then  $\mathcal{M}(d)$  is periodic with period

$$P = \operatorname{lcm} \left\{ \frac{2\pi}{\omega_k} \right\}_{k=1}^K.$$

Moreover, if the  $\omega_k$  are not rational multiples,  $\mathcal{M}(d)$  exhibits quasiperiodic behavior.

The proof of Theorem 1 appears in Appendix A. An optional exponential damping function,  $\mathcal{D}(d) = \exp(-\gamma d)$  with  $\gamma \geq 0$  a learnable coefficient, attenuates attention between distant tokens (Figure 1). Larger  $\gamma$  values promote localized interactions, while smaller ones permit attention across longer ranges, enhancing training stability. The relationship between damping and gradient flow is detailed in Appendix C.

**Theorem 2** (Boundedness and Convergence of Modulated Attention). Let  $S_{ij}$  be the attention score between tokens i and j in FourierRoFormer,

$$\mathbf{S}_{ij} = \frac{\langle \mathbf{q}_i^{RoPE}, \, \mathbf{k}_j^{RoPE} \rangle}{\sqrt{d}} \, \mathcal{M}(d_{ij}) \, e^{-\gamma d_{ij}}, \qquad d_{ij} = |i - j|,$$

where  $\mathcal{M}(d)$  is the Fourier modulation function,  $\gamma>0$  is the damping factor, and  $\|\mathbf{q}_i^{RoPE}\|, \|\mathbf{k}_j^{RoPE}\| \leq M$  for some finite M>0. First, these scores are uniformly bounded, since  $|\mathbf{S}_{ij}| \leq M^2 e^{-\gamma d_{ij}}/\sqrt{d}$ . Second, for any fixed token i, the exponential series of scores converges as the sequence length  $N\to\infty$ , we have  $\sum_{j=1}^N e^{\mathbf{S}_{ij}} < \infty$ . Finally, the corresponding normalized attention weights  $A_{ij} = e^{\mathbf{S}_{ij}}/\sum_{k=1}^N e^{\mathbf{S}_{ik}}$  lie strictly between 0 and 1 for every pair of tokens (i,j), ensuring well-defined probabilistic attention.

Theorem 2 states that attention scores exponentially decrease with distance, so distant tokens have minimal impact on the softmax. Lemma 1 in Appendix B provides additional technical results on the effective attention range. Fourier modulation and damping create a flexible yet structured attention pattern that adapts to visual data while maintaining interpretability, as shown in Figure 1 (right panel). Appendix B shows how FourierRoFormer balances local and global dependencies through Local-Global Balance (Corollary 1), where high-frequency components capture local patterns while low-frequency ones preserve global context. The theoretical analysis demonstrates that the gradient of attention scores with respect to modulation parameters decays exponentially with token distance, ensuring stable training dynamics. Detailed gradient bounds and stability analysis are in Appendix C. Fourier modulation maintains RoPE's geometric properties—translation equivariance, purely relative dependence, and multiplicative separability—within the combined attention mechanism. The proof is in Appendix D.

**Theorem 3** (RoPE-Fourier Compatibility). *In FourierRoFormer, the modulated RoPE attention score* 

$$\mathbf{S}_{mn} = \frac{\langle \mathbf{R}_{\theta,m} \mathbf{q}_m, \mathbf{R}_{\theta,n} \mathbf{k}_n \rangle}{\sqrt{d}} \cdot \mathcal{M}(|m-n|) \cdot e^{-\gamma|m-n|}$$

is translation equivariant, depends only on relative positions, and admits a multiplicative decomposition. Specifically, for any shift  $\tau \in \mathbb{Z}$ , we have  $\mathbf{S}_{(m+\tau)(n+\tau)} = \mathbf{S}_{mn}$ , and  $\mathbf{S}_{mn}$  can be expressed as  $\mathbf{S}_{mn} = f(m-n,\mathbf{q}_m,\mathbf{k}_n)$  for some function f independent of absolute positions. Moreover, the score factorizes as  $\mathbf{S}_{mn} = \mathbf{S}_{mn}^{RoPE} \cdot \mathbf{S}_{mn}^{Fourier}$ , where  $\mathbf{S}_{mn}^{RoPE}$  is the standard RoPE attention score and  $\mathbf{S}_{mn}^{Fourier} = \mathcal{M}(|m-n|) \cdot e^{-\gamma|m-n|}$ .

FourierRoFormer follows the standard Vision Transformer (ViT) pipeline. An input image is first split into fixed-size patches—typically  $4\times 4$  or  $16\times 16$  pixels—which are flattened and linearly projected into the model's embedding dimension. A learnable class (CLS) token is then prepended, and absolute positional embeddings are added to every token in the sequence. The resulting token stream is processed by a stack of Transformer encoder layers. Each layer replaces the usual multihead self-attention with our FourierRoFormer attention module, augments it with rotary positional encoding, and couples it to a feed-forward network, all wrapped in residual connections and layer normalization. After the final encoder block, the representation of the CLS token is fed to a linear classifier that outputs class probabilities.

To isolate the effect of the attention mechanism, we keep every other architectural detail fixed when comparing against baseline models. The vanilla ViT uses standard self-attention with absolute position encodings; DeiT adopts the same backbone while adding a distillation token and

Table 1: ImageNet-1K classification results. FourierRoFormer shows consistent gains across model scales and offers a better performance-parameter trade-off than spectral methods.

Method	Params (M)	<b>GFLOPs</b>	Top-1 (%)	Top-5 (%)				
Standard Vision Transformers								
ViT-B	86.6	17.6	81.8	95.8				
DeiT-B	86.6	17.6	81.8	95.6				
RoFormer-S	22.01	4.60	78.9	94.2				
RoFormer-M	24.75	4.60	81.9	95.7				
RoFormer-B	86.4	17.5	82.3	95.9				
Spectral Methods								
GFNet-H-B	54.0	8.6	82.9	96.1				
WaveViT-B	33.5	6.8	84.8	97.2				
SpectFormer-H-B	33.1	6.3	85.1	97.4				
SVT-H-B	32.8	6.5	85.2	97.3				
FourierRoFormer (Ours)								
FourierRoFormer-S	22.01	4.61	80.4	95.1				
FourierRoFormer-M	24.76	4.63	83.4	96.5				
FourierRoFormer-B	86.41	17.53	84.1	96.9				

teacher–student training; RoFormer swaps in rotary positional embeddings; and FourierRoFormer further enriches RoFormer by superimposing a learnable Fourier modulation and an optional exponential damping term. As shown in Appendix G, this modification preserves the asymptotic computational complexity of the original Transformer and introduces only a modest number of additional parameters, yet it substantially increases representational flexibility.

#### 4 EXPERIMENTAL EVALUATION

We conducted extensive experiments to evaluate FourierRoFormer across multiple scales and tasks, including image classification on CIFAR, ImageNet, object detection and segmentation on COCO, and detailed analysis of learned frequency patterns. Our evaluation is designed to validate both the performance benefits and the theoretical insights developed in the methodology section.

**Experimental Setup.** To evaluate FourierRoFormer's capabilities, we evaluated classification tasks (CIFAR-10/100, ImageNet-1K, Oxford-Flowers102) and dense prediction tasks (object detection, segmentation in COCO) using identical training protocols. We report the mean accuracy on 5 random seeds with statistical significance testing (p < 0.05). For small-scale datasets, we use  $4 \times 4$  patches; ImageNet and COCO use  $16 \times 16$  patches. We evaluated three model sizes: *small* (192d, 6h, 6l), *medium* (384d, 6h, 12l), and *large* (576d, 12h, 12l). FourierRoFormer is initialized with four learnable Fourier components, frequencies linearly spaced between 0.1 and 2.0, amplitudes of 0.1, zero phase, and damping coefficient  $\gamma = 0.01$ .

ImageNet-1K Results. Table 1 presents our comprehensive ImageNet-1K evaluation, including comparisons with spectral methods and scaling analysis. The results demonstrate scaling consistency with improvements ranging from +1.4pp to +1.8pp across model sizes. FourierRoFormer achieves parameter efficiency by using 25% fewer parameters than SpectFormer (24.76M vs 33.1M) while maintaining competitive performance (84.1% vs 85.1%). The approach also provides computational efficiency with 27% lower FLOPs than leading spectral methods. All improvements show statistical significance with p < 0.01 across 5 random seeds.

**Small-Scale Dataset Results.** Table 2 presents comprehensive results on CIFAR and Oxford-Flowers102 in multiple model sizes. The greatest improvements occur in CIFAR-100 (+5.84pp over RoFormer), demonstrating the value of frequency awareness for fine-grained classification tasks with many classes. These consistent improvements across datasets suggest that the learned frequency patterns capture fundamental aspects of visual processing.

**Model Size Scaling Analysis.** To understand how our frequency-aware attention scales with model capacity, Table 3 analyzes performance across different model sizes on CIFAR-100. Notably, our medium-sized FourierRoFormer (84.26%) surpasses even large-sized ViT (81.54%) and DeiT (82.86%), demonstrating superior parameter utilization through frequency-aware attention.

Table 2: Classification results on small-scale datasets. Numbers show mean  $\pm$  standard deviation over 5 independent runs.

Model	CIFAR-10	CIFAR-100	Oxford-Flowers102
Standard ViT DeiT RoFormer	$93.21 \pm 0.14$ $94.58 \pm 0.12$ $94.63 \pm 0.11$	$77.79 \pm 0.21$ $79.55 \pm 0.18$ $78.42 \pm 0.19$	$93.68 \pm 0.18$ $94.75 \pm 0.15$ $94.23 \pm 0.16$
FourierRoFormer	$96.28 \pm 0.10$	$84.26 \pm 0.15$	$\frac{94.23 \pm 0.10}{96.04 \pm 0.13}$

Table 3: Top-1 accuracy on CIFAR-100 across model sizes showing consistent improvements and parameter efficiency.

Model	Small (192d, 6h, 6l)	Medium (384d, 6h, 12l)	Large (576d, 12h, 12l)	Avg Improvement
ViT	$73.62 \pm 0.25$	$77.79 \pm 0.21$	$81.54 \pm 0.17$	-
DeiT	$75.28 \pm 0.23$	$79.55 \pm 0.18$	$82.86 \pm 0.16$	-
RoFormer	$76.04 \pm 0.22$	$78.42 \pm 0.19$	$82.97 \pm 0.15$	-
FourierRoFormer	$80.39 \pm 0.19$	$84.26\pm0.15$	$86.52 \pm 0.13$	+4.8pp
Improvement	+4.35pp	+5.84pp	+3.55pp	-

**Object Detection and Segmentation Results.** We evaluate on COCO using Mask R-CNN with FourierRoFormer as the backbone, expecting larger improvements due to the multi-scale nature of detection tasks (see Table 4). The largest improvements occur on medium-scale objects (+5.1pp) where frequency awareness provides maximum benefit, confirming our hypothesis about multi-scale reasoning advantages.

## Comprehensive Ablation Studies.

Component Analysis. Our ablations reveal that Fourier modulation provides larger benefit (+4.43pp) than damping (+2.09pp), yet the components work complementarily to achieve +5.84pp total improvement over the baseline. The optimal configuration uses 4-8 Fourier components with moderate damping ( $\gamma = 0.01$ ). See Appendix F for detailed analysis in Table 16.

Frequency Initialization Strategies. Among various initialization approaches, logarithmic spacing achieves best performance (+0.36pp over linear spacing), providing better coverage of the frequency spectrum. Complete results are presented in Appendix F in Table 17.

Multi-Head Frequency Specialization Analysis. One of our key findings is that different attention heads learn distance-based attention patterns when given independent parameters. To analyze the relationship between learned frequencies and visual patterns, we compute attention maps for 1,000 randomly sampled validation images. For each attention head, we: (1) extract the dominant frequency component based on amplitude, (2) segment images using ground-truth object masks when available or edge detection (Canny) otherwise, (3) compute Pearson correlation between attention weights and binary masks for boundaries/textures/global regions. The reported correlations represent averages across the validation sample.

Our analysis shows that heads 1-2 predominantly use low frequencies (0.2-0.6 Hz) with attention spanning approximately 89 tokens, while heads 3-4 employ mid frequencies (0.6-1.4 Hz) with attention focused on approximately 43 tokens. Finally, heads 5-6 utilize high frequencies (1.4-3.2 Hz) to handle fine details within 21 tokens. This specialization emerges after 35 epochs and stabilizes by epoch 100, providing evidence of learned frequency-based division of labor. These findings suggest that the model automatically discovers an optimal division of attention across different spatial scales. Complete results are presented in Appendix F in Table 15.

Table 4: COCO object detection and instance segmentation results showing FourierRoFormer's advantages for multi-scale tasks.

Backbone	Detection mAP	Segmentation mAP	Medium Objects	Small Objects
RoFormer FourierRoFormer	41.2 <b>43.8</b>	37.9 <b>40.1</b>	22.4 <b>27.5</b>	15.8 <b>18.9</b>
Improvement	+2.6pp	+2.2pp	+5.1pp	+3.1pp

Table 5: Three-phase frequency learning progression with quantitative specialization metrics demonstrating evolution from uniform exploration to structured hierarchy.

Phase	Epochs	Coeff. Var.	Entropy	Stability	Freq Variance	Corr.	Convergence
Exploration	0-40	0.12	$3.41 \pm 0.18$	< 30%	0.08	0.34	Unstable
Specialization	40-120	0.68	$3.38 \pm 0.12$	70%	0.31	0.67	Progressing
Convergence	120+	0.91	$3.35 \pm 0.08$	> 95%	0.42	0.84	Stable

Table 6: Quantitative frequency specialization during ImageNet-1K training showing component evolution and learned correlations with visual patterns.

Component	Initial Amp	Final Amp	Learned Freq	Visual Pattern	Correlation
k=1	$0.10 \pm 0.02$	0.43	0.3 Hz	Global shape	r = 0.78
k=2	$0.10 \pm 0.02$	0.31	1.1 Hz	Object boundaries	r = 0.85
k=3	$0.10 \pm 0.02$	0.18	2.4 Hz	Fine textures	r = 0.71
k=4	$0.10 \pm 0.02$	0.08	3.2 Hz	Noise/artifacts	r = 0.34

**Training Dynamics and Frequency Learning Validation.** To further validate our theoretical predictions about frequency learning, we provide concrete empirical evidence through detailed analysis of training dynamics and component evolution.

We systematically tracked all Fourier component parameters at 10-epoch intervals throughout training across all 5 runs. The phases identified represent consistent patterns observed across runs, not post-hoc categorization. Specifically, we measure the coefficient of variation (CV) of amplitudes, the parameter update magnitude via  $\ell_2$  norm, and the attention entropy. Phase boundaries are defined by thresholds:

Exploration: CV < 0.3, Specialization:  $0.3 \le CV < 0.7$ , Convergence:  $CV \ge 0.7$ .

- Phase 1 (Epochs 0-40): Exploration. This phase exhibits uniform amplitude distribution with coefficient of variation = 0.12 indicating low specialization, where all frequency components contribute equally ( 25% each). The model shows high variance in attention patterns (entropy =  $3.41 \pm 0.18$ ) and explores different spatial scales with < 30% parameter stability. Pattern correlation remains weak (r = 0.34) indicating random exploration.
- Phase 2 (Epochs 40-120): Specialization. This phase shows emerging specialization as the coefficient of variation increases to 0.68 indicating moderate specialization. A clear frequency hierarchy emerges with 70% parameter stability, while attention becomes more structured (entropy =  $3.38 \pm 0.12$ ). Components begin correlating with visual patterns (r = 0.67) and frequency variance increases to 0.31 showing differentiation.
- Phase 3 (Epochs 120+): Convergence. The final phase demonstrates strong specialization with coefficient of variation = 0.91 indicating near-maximal specialization and stable frequency allocation with > 95% parameter stability. The model achieves strong pattern-frequency correlation (r = 0.84) indicating semantic alignment, optimal attention structure (entropy =  $3.35 \pm 0.08$ ), and maximum frequency variance (0.42) showing complete differentiation.

Specialzation Metrics Definition: Coefficient of variation measures amplitude dispersion ( $CV=\sigma/\mu$ ), where higher values indicate stronger component differentiation. Stability percentage tracks parameter convergence, and pattern correlation measures alignment with ground-truth visual patterns.

Quantitative Frequency Component Analysis. Table 6 shows how different components specialize during training. This quantitative analysis confirms that different frequency components learn to capture complementary visual patterns, with the strongest correlation (r = 0.85) for object boundary detection at 1.1 Hz.

Comprehensive Efficiency Analysis. Table 7 provides detailed computational efficiency comparison. We use  $efficiency\ Score = \frac{\text{Top-1 Accuracy}}{\log(\text{Params}) \times \sqrt{\text{Training Time}}}$  that captures performance-complexity trade-

offs. The analysis reveals improved parameter efficiency with only 0.04% parameter overhead for 1.5pp accuracy gain. The approach is memory efficient with minimal memory increase (0.6%) compared to parameter gains and maintains training stability with similar training time but improved convergence. Overall, it achieves a superior tradeoff with 17% better efficiency score than Ro-Former baseline. Table 8 compares FourierRoFormer with recent advances in positional encoding.

Table 7: Comprehensive efficiency analysis showing FourierRoFormer's minimal overhead for gains.

Method	Params (M)	Memory (GB)	Throughput (img/s)	Training Time (h)	Top-1 (%)	Efficiency Score
RoFormer-M	24.75	18.0	220	12.0	81.9	3.33
GFNet-H-B	54.0	21.5	185	16.8	82.9	2.41
SpectFormer-H-B	33.1	19.2	195	14.5	85.1	3.21
FourierRoFormer-M	24.76	18.1	215	12.3	83.4	3.91
Overhead vs RoFormer	+0.04%	+0.6%	-2.3%	+2.5%	+1.5pp	+17%

FourierRoFormer's key advantage is learning adaptive frequency patterns rather than using fixed biases or interpolation schemes, which leads us to investigate the underlying mechanisms that drive these performance gains.

Table 8: Comparison with recent positional encoding methods on ImageNet-1K showing advantages of learnable frequency patterns.

Method	Description	Top-1 (%)	Key Limitation
ALiBi Context-aware Biases Functional Interpolation RoFormer	Linear bias attention Length extrapolation focus RoPE interpolation Rotary embeddings	82.7 83.1 83.4 82.3	Fixed linear decay Limited frequency awareness No adaptive patterns Uniform frequency treatment
FourierRoFormer	Learnable frequency patterns	84.1	-
Improvement	vs best baseline	+0.7pp	Adaptive learning

#### 5 Analysis and Discussion

Having established FourierRoFormer's empirical advantages, we now turn to understanding the mechanisms behind these improvements and analyzing how the model learns to leverage frequency information.

**Frequency Learning Mechanism Understanding.** Our approach enables the model to learn optimal spatial frequencies that align with natural image statistics.

- Adaptive Scale Discovery. The model automatically discovers optimal spatial frequencies (0.3, 1.1, 2.4 Hz) that correspond to different visual scales global structure, object boundaries, and fine details respectively.
- Attention Structure Emergence. Learned frequencies create attention patterns that strongly correlate (r = 0.85) with ground-truth object boundaries, demonstrating semantic alignment:
- Hierarchical Processing. Low frequencies (0.3 Hz) capture global context across 89 tokens, while high frequencies (2.4 Hz) focus on local details within 21 tokens, creating natural hierarchical attention.

**Post-Attention Modulation Design Justification.** Our choice to apply Fourier modulation after attention computation (rather than before) is theoretically and empirically motivated:

- <u>Theoretical Justification</u>. Post-attention modulation preserves semantic query-key relationships while adding frequency awareness. Pre-attention modulation disrupts learned embedding geometry that encodes semantic similarity.
- Empirical Evidence. Experiments show post-attention achieves 84.1% vs 82.3% for pre-attention (- $\overline{1.8}$ pp), with more stable gradients ( $\sigma = 0.12$  vs 0.41 for pre-attention). Post-attention maintains stable gradient magnitudes across layers, while pre-attention causes 34% higher gradient variance.

Comparison with Spectral Transformer Methods. Our approach offers distinct advantages over existing spectral methods (see Table 9). It is a learnable adaptation by learning data-specific frequency patterns, unlike fixed Fourier (GFNet) or wavelet (WaveViT) transforms. It maintains architectural simplicity by preserving the standard transformer structure unlike methods that require architectural overhaul. The method provides a strong theoretical foundation with formal guarantees for boundedness, convergence, and interpretability, while achieving competitive accuracy with significantly fewer parameters.

Table 9: Detailed comparison with spectral transformer methods showing FourierRoFormer's unique advantages.

Feature	GFNet	WaveViT	SpectFormer	SVT	FourierRoFormer
Adaptive frequency selection	Х	<b>√</b> (wavelet)	<b>√</b> (limited)	<b>√</b> (wavelet)	<b>√</b> (learned)
Interpretable modulation	X	X	×	×	✓
Learnable damping & stability	X	×	X	×	✓
Theoretical guarantees	X	X	×	×	✓
Architecture compatibility	X	Moderate	Moderate	×	✓
Parameter efficiency	Moderate	Moderate	Good	Good	Excellent

**Attention Pattern Visualization and Analysis.** Our visualizations reveal that FourierRoFormer produces highly structured attention patterns that align with semantic image content:

- Standard ViT. Produces diffused, weakly structured attention with limited semantic alignment.
- <u>RoFormer</u>. Shows improved spatial awareness through relative position encoding but still covers broad, unfocused regions.
- <u>FourierRoFormer</u>. Exhibits highly structured attention emphasizing object boundaries and semantic features, with clear multi-scale organization.

The frequency-aware modulation creates natural attention hierarchies where different frequency components focus on complementary spatial scales, resulting in more interpretable and effective visual processing.

**Implications for Transformer Design.** These findings have broader implications for transformer architecture design. The success of learned frequency modulation suggests that incorporating domain-specific inductive biases through principled mathematical frameworks can significantly enhance model performance while maintaining interpretability. The approach bridges data-driven learning with structured frequency-based inductive biases, offering a principled way to embed multiscale spatial awareness in transformers.

## 6 Conclusion

We introduced FourierRoFormer, a transformer architecture that incorporates learnable Fourier components to bring frequency awareness to the attention mechanism. Our approach enables the adaptive capture of multi-scale visual patterns while maintaining architectural simplicity and theoretical rigor. Comprehensive evaluations demonstrate consistent improvements, with our model achieving 84.1% top-1 accuracy on ImageNet-1K (+1.8pp over RoFormer-B) and significant gains on fine-grained classification and dense prediction tasks.

Our key contributions are: (1) a novel mechanism for learning adaptive frequency patterns directly within attention scores; (2) theoretical guarantees for the model's expressivity, stability, and interpretability; and (3) a detailed analysis revealing that the model learns a functional hierarchy of frequencies, where different attention heads specialize in distinct spatial scales.

While the base model already performs strongly, our analysis shows that allowing head-specific frequency parameters yields further accuracy gains (+0.5pp), confirming the value of specialization. Key limitations include the inherited  $\mathcal{O}(n^2)$  complexity of standard attention. Future work will focus on integrating these head-specific patterns into sparse attention variants and extending the frequency-aware framework to other modalities like video, speech, and language. Ultimately, FourierRoFormer bridges data-driven learning with principled, frequency-based inductive biases, offering a robust method for embedding multi-scale awareness in transformers.

#### REFERENCES

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

Krzysztof Choromanski, Valerii Likhosherstov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention with performers. *arXiv preprint arXiv:2009.14794*, 2020.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

Yanyu Li, Geng Yuan, Yang Wen, Ju Hu, Georgios Evangelidis, Sergey Tulyakov, Yanzhi Wang, and Jian Ren. Efficientformer: Vision transformers at mobilenet speed. *Advances in Neural Information Processing Systems*, 35:12934–12949, 2022.

- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
- Sachin Mehta and Mohammad Rastegari. Mobilevit: light-weight, general-purpose, and mobile-friendly vision transformer. *arXiv* preprint arXiv:2110.02178, 2021.
- Edouard Oyallon, Eugene Belilovsky, Sergey Zagoruyko, and Michal Valko. Compressing the input for cnns with the first-order scattering transform. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 301–316, 2018.
- Namuk Park and Songkuk Kim. How do vision transformers work? In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=D78Go4hVcxO.
- Badri Patro and Vijay Agneeswaran. Scattering vision transformer: Spectral mixing matters. *Advances in Neural Information Processing Systems*, 36:54152–54166, 2023.
- Badri N Patro, Vinay P Namboodiri, and Vijay S Agneeswaran. Spectformer: Frequency and attention is what you need in a vision transformer. In 2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), pp. 9543–9554. IEEE, 2025a.
- Badri N. Patro, Vinay P. Namboodiri, and Vijay S. Agneeswaran. Spectformer: Frequency and attention is what you need in a vision transformer. In 2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), pp. 9543–9554, 2025b. doi:10.1109/WACV61041.2025.00924.
- Allan Pinkus. Approximation theory of the mlp model in neural networks. *Acta Numerica*, 8: 143–195, 1999. doi:10.1017/S0962492900002937.
- Ofir Press, Noah A Smith, and Mike Lewis. Train short, test long: Attention with linear biases enables input length extrapolation. *arXiv preprint arXiv:2108.12409*, 2021.
- Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? *Advances in neural information processing systems*, 34:12116–12128, 2021.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In John C. Platt, Daphne Koller, Yoram Singer, and Sam T. Roweis (eds.), *Advances in Neural Information Processing Systems*, volume 20, pp. 1177–1184. Curran Associates, Inc., 2008. URL https://papers.nips.cc/paper\_files/paper/2007/file/013a006f03dbc5392effeb8f18fda755-Paper.pdf.
- Yongming Rao, Wenliang Zhao, Zheng Zhu, Jiwen Lu, and Jie Zhou. Global filter networks for image classification. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- Walter Rudin. Principles of Mathematical Analysis. McGraw-Hill, New York, 3rd edition, 1976.
- Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. *arXiv preprint arXiv:1803.02155*, 2018.
  - Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.

Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Ren Ng, and Jonathan T. Barron. Fourier features let networks learn high frequency functions in low dimensional domains. In *Advances in Neural Information Processing Systems*, volume 33, pp. 7537–7547. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper/2020/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.

Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In International conference on machine learning, pp. 10347–10357. PMLR, 2021.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity, 2020.

Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 568–578, 2021.

Ting Yao, Yingwei Pan, Yehao Li, Chong-Wah Ngo, and Tao Mei. Wave-vit: Unifying wavelet and transformers for visual representation learning. In Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner (eds.), *Computer Vision – ECCV 2022*, pp. 328–345, Cham, 2022. Springer Nature Switzerland. ISBN 978-3-031-19806-9.

## A ANALYSIS OF FOURIER MODULATION FUNCTION

The FourierRoFormer introduces a learned mixture of sinusoidal components to modulate attention based on token distances. We first analyze the properties of this modulation function and establish its theoretical guarantees.

**Theorem 1** (Properties of Fourier Modulation Function). Let  $\mathcal{M}:\mathbb{R}\to(0,1)$  be the Fourier modulation function defined as

$$\mathcal{M}(d) = \frac{1}{2} \left( \tanh \left( \sum_{k=1}^{K} a_k \cos(\omega_k d + \phi_k) \right) + 1 \right)$$

where  $a_k \in \mathbb{R}$  are learnable amplitudes,  $\omega_k > 0$  are learnable frequencies, and  $\phi_k \in [0, 2\pi)$  are learnable phase shifts for  $k = 1, \ldots, K$ . Then  $\mathcal{M}(d)$  is a smooth function with  $\mathcal{M}(d) \in (0, 1)$  for all  $d \in \mathbb{R}$ . For any continuous function  $f : [0, L] \to (0, 1)$  and any  $\varepsilon > 0$ , there exists an integer K and parameters  $\{a_k, \omega_k, \phi_k\}_{k=1}^K$  such that

$$\sup_{d \in [0,L]} |\mathcal{M}(d) - f(d)| < \varepsilon$$

If the set of frequencies  $\{\omega_k\}_{k=1}^K$  consists of rational multiples of each other, then  $\mathcal{M}(d)$  is periodic with period

$$P = \operatorname{lcm} \left\{ \frac{2\pi}{\omega_k} \right\}_{k=1}^K$$

Moreover, if the  $\omega_k$  are not rational multiples,  $\mathcal{M}(d)$  exhibits quasiperiodic behavior.

*Proof.* We prove each part in turn. For any  $x \in \mathbb{R}$ , it holds that  $\tanh(x) \in (-1,1)$ . Consider the inner sum:

$$S(d) = \sum_{k=1}^{K} a_k \cos(\omega_k d + \phi_k)$$

Since  $cos(\theta) \in [-1, 1]$  for all  $\theta \in \mathbb{R}$ , we have:

$$|S(d)| \le \sum_{k=1}^{K} |a_k|$$

Thus,  $\tanh(S(d)) \in (-1,1)$  for all  $d \in \mathbb{R}$ . Applying the affine transformation  $x \mapsto \frac{1}{2}x + \frac{1}{2}$  maps (-1,1) to (0,1):

$$\mathcal{M}(d) = \frac{1}{2} \left( \tanh(S(d)) + 1 \right) \in (0, 1)$$

Furthermore, since  $\cos$ ,  $\tanh$ , and affine transformations are smooth functions,  $\mathcal{M}(d)$  is infinitely differentiable, i.e.,  $\mathcal{M} \in C^{\infty}(\mathbb{R})$ . Let  $f: [0, L] \to (0, 1)$  be continuous. Define the lifted function:

$$g(d) = \tanh^{-1} \left( 2f(d) - 1 \right)$$

Note that since  $f(d) \in (0,1)$ , we have  $2f(d)-1 \in (-1,1)$ , and thus g(d) is well-defined and continuous on [0,L]. By the Stone–Weierstrass theorem, the algebra of trigonometric polynomials is dense in the space of continuous real-valued functions on [0,L] (see, e.g., Rudin (1976)). Moreover, the use of nonlinear activation functions applied to sinusoidal expansions falls within the scope of classical approximation theory for neural networks Pinkus (1999). Therefore, for any  $\varepsilon'>0$ , there exist parameters  $\{a_k,\omega_k,\phi_k\}_{k=1}^K$  such that

$$\sup_{d \in [0,L]} \left| g(d) - \sum_{k=1}^{K} a_k \cos(\omega_k d + \phi_k) \right| < \varepsilon'$$

Since  $\tanh$  is continuous and Lipschitz on compact sets, there exists a constant  $L_{\tanh}$  such that:  $|\tanh(x) - \tanh(y)| \le L_{\tanh}|x-y|$  for all x,y in the image of g(d) and its approximation. Thus, we have:

$$\sup_{d \in [0,L]} \left| \tanh(g(d)) - \tanh\left(\sum_{k=1}^{K} a_k \cos(\omega_k d + \phi_k)\right) \right| < L_{\tanh} \varepsilon'$$

Multiplying by  $\frac{1}{2}$  and adding  $\frac{1}{2}$  preserves the approximation margin. By choosing  $\varepsilon' = \frac{\varepsilon}{L_{\rm tanh}}$ , we ensure:

$$\sup_{d \in [0,L]} |f(d) - \mathcal{M}(d)| < \varepsilon$$

Thus,  $\mathcal{M}(d)$  uniformly approximates any continuous function f on [0,L] to arbitrary precision. Each term  $\cos(\omega_k d + \phi_k)$  is periodic with period  $\frac{2\pi}{\omega_k}$ . If all frequencies  $\omega_k$  are rational multiples of each other, there exists a common period:

$$P = \operatorname{lcm} \left\{ \frac{2\pi}{\omega_k} \right\}_{k=1}^K$$

Thus, the finite sum S(d) is periodic with period P. Since  $\tanh$  and affine transformations are applied pointwise and preserve periodicity,  $\mathcal{M}(d)$  is also periodic with period P.

In addition to the approximation and periodicity properties established above, the form of  $\mathcal{M}(d)$  provides clear interpretability of the roles played by its parameters, as summarized in the following corollary.

**Proposition 2** (Interpretability of Fourier Components). The learned parameters  $\{a_k, \omega_k, \phi_k\}_{k=1}^K$  in the modulation function  $\mathcal{M}(d)$  admit the following interpretations:

- Amplitude( $a_k$ ) controls the contribution strength of the k-th frequency component to the overall modulation pattern. Larger  $|a_k|$  values amplify the influence of the corresponding cosine term.
- Frequency ( $\omega_k$ ) determines the spatial frequency of the oscillations, i.e., how rapidly the attention modulation varies with respect to token distance d. Higher  $\omega_k$  yields finer-grained, higher-frequency patterns.
- Phase shift  $(\phi_k)$  specifies the horizontal displacement of the k-th component along the distance axis, enabling translation of attention peaks and troughs without altering their frequency.

The interpretability of  $\{a_k, \omega_k, \phi_k\}_{k=1}^K$  facilitates analysis of learned attention patterns and enables explicit control over the modulation behavior. For example, sparsity-promoting regularization on  $\{a_k\}$  can encourage parsimonious attention structures.

*Proof.* We examine the modulation function:

$$\mathcal{M}(d) = \frac{1}{2} \left( \tanh \left( \sum_{k=1}^{K} a_k \cos(\omega_k d + \phi_k) \right) + 1 \right)$$

and analyze the role of each parameter  $\{a_k, \omega_k, \phi_k\}$  in shaping  $\mathcal{M}(d)$ . Consider the inner argument of the tanh function:

$$S(d) = \sum_{k=1}^{K} a_k \cos(\omega_k d + \phi_k)$$

This is a finite sum of cosine functions, each parameterized by amplitude, frequency, and phase shift. The amplitude  $a_k$  scales the contribution of the k-th component: increasing  $|a_k|$  amplifies its oscillatory magnitude, while the sign determines whether it reinforces or counteracts other terms. The frequency  $\omega_k$  controls the spatial scale, with the component completing one full oscillation over  $T_k = \frac{2\pi}{\omega_k}$ ; larger  $\omega_k$  produces finer, more rapid oscillations over token distance d. The phase shift  $\phi_k$  translates the cosine along the d-axis, corresponding to a horizontal displacement of  $\Delta d = -\phi_k/\omega_k$ , which adjusts the positions of peaks and troughs without affecting amplitude or frequency.

Finally, observe that the outer tanh function is a smooth, monotonically increasing function applied pointwise to S(d). While tanh compresses the range of S(d) into (-1,1), it preserves the relative locations of maxima, minima, and zero crossings of S(d), thereby maintaining the interpretability of the underlying sinusoidal components. The subsequent affine transformation maps this range to (0,1) without altering these relationships. Thus, the parameters  $\{a_k,\omega_k,\phi_k\}_{k=1}^K$  maintain clear and interpretable roles in controlling the shape and characteristics of  $\mathcal{M}(d)$ .

## B CONVERGENCE ANALYSIS OF MODULATED ATTENTION

We now analyze how the Fourier modulation influences attention scores and their convergence behavior, particularly focusing on the boundedness of scores, the normalization of attention weights, and their behavior as the sequence length grows. The following theorem establishes uniform bounds and guarantees well-posedness of the attention mechanism in FourierRoFormer.

**Theorem 2** (Boundedness and Convergence of Modulated Attention). Let  $S_{ij}$  denote the attention score between tokens i and j in FourierRoFormer, defined as

$$\mathbf{S}_{ij} = \frac{\langle \mathbf{q}_i^{RoPE}, \mathbf{k}_j^{RoPE} \rangle}{\sqrt{d}} \cdot \mathcal{M}(d_{ij}) \cdot e^{-\gamma d_{ij}}$$

where  $d_{ij} = |i - j|$ ,  $\mathcal{M}(d)$  is the Fourier modulation function,  $\gamma > 0$  is the damping factor, and  $\|\mathbf{q}_i^{RoPE}\|$ ,  $\|\mathbf{k}_j^{RoPE}\| \le M$  for some finite constant M > 0. Then, the following properties hold:

1. The attention scores are bounded:

$$|\mathbf{S}_{ij}| \le \frac{M^2}{\sqrt{d}} \, e^{-\gamma d_{ij}}$$

2. For any fixed token i, as sequence length  $N \to \infty$ ,

$$\sum_{j=1}^{N} e^{\mathbf{S}_{ij}} < \infty$$

3. For all pairs (i, j), the normalized attention satisfies

$$A_{ij} = \frac{e^{\mathbf{S}_{ij}}}{\sum_{k=1}^{N} e^{\mathbf{S}_{ik}}} \in (0,1).$$

*Proof.* We prove each part in turn. First, by the Cauchy–Schwarz inequality, and under the assumption  $\|\mathbf{q}_i^{\text{RoPE}}\|$ ,  $\|\mathbf{k}_i^{\text{RoPE}}\| \le M$ , we have:

$$|\langle \mathbf{q}_i^{\text{RoPE}}, \mathbf{k}_j^{\text{RoPE}} \rangle| \le M^2$$

From Theorem 1,  $\mathcal{M}(d_{ij}) \in (0,1)$  for all  $d_{ij}$ , and by definition, the damping factor is  $\mathcal{D}(d_{ij}) = e^{-\gamma d_{ij}}$ . Hence:

$$|\mathbf{S}_{ij}| \leq \frac{M^2}{\sqrt{d}} e^{-\gamma d_{ij}}$$

To show the convergence of the normalization sum, we use the below estimate:

$$\sum_{j=1}^{N} e^{\mathbf{S}_{ij}} \le \sum_{j=1}^{N} \exp\left(\frac{M^2}{\sqrt{d}} e^{-\gamma|i-j|}\right)$$

Since  $e^{-\gamma|i-j|} \to 0$  exponentially as  $|i-j| \to \infty$ , and  $\exp\left(c\,e^{-\gamma|i-j|}\right) \to 1$ , the summand behaves like a constant for small |i-j| and decays exponentially for large |i-j|. Thus, the sum can be split:

$$\sum_{j \le i} \exp\left(\frac{M^2}{\sqrt{d}} e^{-\gamma(i-j)}\right) + \sum_{j > i} \exp\left(\frac{M^2}{\sqrt{d}} e^{-\gamma(j-i)}\right)$$

Each term is a convergent exponential series, as  $e^{-\gamma n}$  decays exponentially and  $\exp\left(c\,e^{-\gamma n}\right)$  remains summable for c>0. This follows from standard results on the convergence of rapidly decreasing exponential series (Rudin, 1976, p. 5). Therefore, the total sum converges as  $N\to\infty$ . The denominator of the attention weights is strictly positive and finite. Moreover, since the numerator  $e^{\mathbf{S}_{ij}}>0$ , it follows that:

$$A_{ij} = \frac{e^{\mathbf{S}_{ij}}}{\sum_{k=1}^{N} e^{\mathbf{S}_{ik}}} \in (0,1)$$

for all i and j. This ensures that attention weights are well-defined probability distributions over tokens.

Building on the boundedness of attention weights, we now characterize the effective receptive field of FourierRoFormer, showing that attention to distant tokens decays below any desired threshold.

**Lemma 1** (Effective Attention Range). For any  $\epsilon > 0$ , there exists a distance  $R_{\epsilon}$  such that for all  $d_{ij} > R_{\epsilon}$ :

$$A_{ij} < \epsilon$$

where  $R_{\epsilon}$  depends on the model parameters  $\{M, d, \gamma, \{a_k, \omega_k, \phi_k\}_{k=1}^K\}$ .

*Proof.* From the bound in Theorem 2(a):

$$\mathbf{S}_{ij} \leq \frac{M^2}{\sqrt{d}} \cdot \exp(-\gamma d_{ij})$$

The attention weight  $A_{ij}$  is bounded by:

$$A_{ij} \le \frac{\exp(\frac{M^2}{\sqrt{d}} \cdot \exp(-\gamma d_{ij}))}{\exp(\frac{M^2}{\sqrt{d}})} = \exp\left(\frac{M^2}{\sqrt{d}} (\exp(-\gamma d_{ij}) - 1)\right)$$

For any  $\epsilon > 0$ , we can solve:

$$\exp\left(\frac{M^2}{\sqrt{d}}(\exp(-\gamma R_{\epsilon}) - 1)\right) = \epsilon$$

This yields:

$$R_{\epsilon} = -\frac{1}{\gamma} \ln \left( 1 + \frac{\sqrt{d}}{M^2} \ln(\epsilon) \right)$$

For  $d_{ij} > R_{\epsilon}$ , we have  $A_{ij} < \epsilon$  by monotonicity.

 The decomposition of the attention modulation into distinct frequency components, together with exponential damping, enables FourierRoFormer to simultaneously capture both fine-grained local patterns and broad global context, as formalized in the following corollary.

**Corollary 1** (Local-Global Balance). The FourierRoFormer attention mechanism balances local and global dependencies through its modulation design: high-frequency Fourier components capture local patterns, low-frequency components preserve global context, and the exponential damping term  $\exp(-\gamma d_{ij})$  ensures smooth decay of attention with distance.

*Proof.* The result follows from the structure of the attention score  $S_{ij}$ , which combines Fourier modulation and exponential damping. First, the high-frequency components with  $\omega_k \gg 1$  induce rapid oscillations in  $\mathcal{M}(d_{ij})$ , enhancing sensitivity to local variations in token distance. Conversely, low-frequency components with  $\omega_k \approx 1$  produce slowly varying modulation, preserving global contextual information. Additionally, the damping factor  $\exp(-\gamma d_{ij})$  enforces an overall decay of attention scores with distance, ensuring that contributions from distant tokens diminish smoothly. Together, these elements balance fine-grained local interactions and long-range global dependencies, while keeping attention scores bounded.

In summary, Theorems 2, Lemma 1, and Corollary 1 establish that FourierRoFormer's attention is bounded, localized, and balances local and global context via its modulation structure. These properties ensure scalability and stability, especially for long sequences.

## C GRADIENT ANALYSIS

In this section we characterize the gradient behavior of the FourierRoFormer modulation parameters, deriving uniform bounds that govern the learning dynamics and inform convergence properties.

**Proposition 3** (Gradient Bounds for Modulation Parameters). Let  $\theta = \{a_k, \omega_k, \phi_k\}_{k=1}^K$  denote the Fourier modulation parameters, and let  $\mathbf{S}_{ij}$  be the attention score between tokens i and j, associated with distance  $d_{ij}$ . Assume the modulation output is scaled by a constant M > 0, and let  $\gamma > 0$  be the effective decay rate. Then, the following gradient bounds hold for all  $k = 1, \ldots, K$ :

(a) Amplitude gradients

$$\left\| \frac{\partial \mathbf{S}_{ij}}{\partial a_k} \right\| \le \frac{M^2}{2\sqrt{d}} e^{-\gamma d_{ij}}$$

(b) Frequency gradients

$$\left\| \frac{\partial \mathbf{S}_{ij}}{\partial \omega_k} \right\| \le \frac{M^2}{2\sqrt{d}} \cdot d_{ij} \, e^{-\gamma d_{ij}}$$

(c) Phase gradients

$$\left\| \frac{\partial \mathbf{S}_{ij}}{\partial \phi_k} \right\| \le \frac{M^2}{2\sqrt{d}} \, e^{-\gamma d_{ij}}$$

*Proof.* We analyze each gradient component individually.

Let  $S_{ij}$  denote the attention score between tokens i and j, with  $d_{ij}$  their distance. Recall:

$$\mathbf{S}_{ij} = \frac{\langle \mathbf{q}_i^{\text{RoPE}}, \mathbf{k}_j^{\text{RoPE}} \rangle}{\sqrt{d}} \cdot \mathcal{D}(d_{ij}) \cdot \mathcal{M}(d_{ij})$$

where  $\mathcal{D}(d_{ij})$  is a distance-dependent decay term, and  $\mathcal{M}(d_{ij})$  is the Fourier modulation function.

For all cases, we use the bound:

$$\left| \frac{\langle \mathbf{q}_i^{\text{RoPE}}, \mathbf{k}_j^{\text{RoPE}} \rangle}{\sqrt{d}} \cdot \mathcal{D}(d_{ij}) \right| \leq \frac{M^2}{\sqrt{d}} \cdot e^{-\gamma d_{ij}}$$

where M>0 bounds the norm of query and key vectors, and  $\gamma>0$  controls the decay. We compute derivatives of  $\mathcal{M}$ , recalling:

$$\mathcal{M}(d) = \frac{1}{2} \left( \tanh(x) + 1 \right), \quad x = \sum_{l=1}^{K} a_l \cos(\omega_l d + \phi_l)$$

Noting that  $\tanh'(x) = 1 - \tanh^2(x)$ , and  $|\tanh'(x)| \le 1$ , we proceed with the amplitude gradients:

$$\frac{\partial \mathcal{M}}{\partial a_k} = \frac{1}{2} \cdot (1 - \tanh^2(x)) \cdot \cos(\omega_k d + \phi_k)$$

Since  $|\cos(\cdot)| \le 1$ , we have:

$$\left\| \frac{\partial \mathbf{S}_{ij}}{\partial a_k} \right\| \le \frac{M^2}{2\sqrt{d}} \cdot e^{-\gamma d_{ij}}$$

Next we look evaluate the frequency gradients:

$$\frac{\partial \mathcal{M}}{\partial \omega_k} = -\frac{1}{2} \cdot (1 - \tanh^2(x)) \cdot a_k d \sin(\omega_k d + \phi_k)$$

Using  $|\sin(\cdot)| \le 1$ , we obtain:

$$\left\| \frac{\partial \mathbf{S}_{ij}}{\partial \omega_k} \right\| \le \frac{M^2}{2\sqrt{d}} \cdot d_{ij} \cdot e^{-\gamma d_{ij}}$$

Finally we estimate the phase gradients:

$$\frac{\partial \mathcal{M}}{\partial \phi_k} = -\frac{1}{2} \cdot (1 - \tanh^2(x)) \cdot a_k \sin(\omega_k d + \phi_k)$$

Thus,

$$\left\| \frac{\partial \mathbf{S}_{ij}}{\partial \phi_k} \right\| \le \frac{M^2}{2\sqrt{d}} \cdot e^{-\gamma d_{ij}}$$

This completes the proof.

Building on the component-wise gradient bounds established in Theorem 3, we now state a general decay property that holds uniformly for all modulation parameters.

**Lemma 2** (Gradient Decay). The gradients of attention scores with respect to Fourier parameters decay exponentially with token distance:

$$\left\| \frac{\partial \mathbf{S}_{ij}}{\partial \theta} \right\| \le C_{\theta} \cdot \exp(-\gamma d_{ij})$$

where  $C_{\theta}$  is a constant depending on the parameter type  $\theta \in \{a_k, \omega_k, \phi_k\}$ .

*Proof.* The result follows directly from Theorem 3. For amplitude and phase parameters, we set  $C_{\theta} = \frac{M^2}{2\sqrt{d}}$ . For frequency parameters, observe that the term  $d_{ij} \cdot e^{-\gamma d_{ij}}$  attains its maximum at  $d_{ij} = 1/\gamma$ , giving  $C_{\theta} = \frac{M^2}{2\gamma e \sqrt{d}}$ .

The exponential gradient decay established in Lemma 2 directly implies desirable properties for the learning dynamics of FourierRoFormer, summarized in the following corollary.

Corollary 2 (Training Stability). Under the exponential gradient decay established in Lemma 2, the training dynamics of FourierRoFormer exhibit the following properties: the magnitude of parameter updates remains bounded throughout training, ensuring stability. The impact of distant tokens on parameter gradients diminishes exponentially with token distance, promoting localized learning. Backpropagation through attention layers remains well-conditioned, preventing gradient explosion or vanishing.

*Proof.* By Lemma 2, the gradient of the attention score with respect to any Fourier parameter  $\theta$  satisfies

 $\left\| \frac{\partial \mathbf{S}_{ij}}{\partial \theta} \right\| \le C_{\theta} \cdot e^{-\gamma d_{ij}}$ 

for some constant  $C_{\theta} > 0$ .

Summing over all token pairs (i, j), the total gradient norm satisfies:

$$\|\nabla_{\theta} \mathcal{L}\| \le C_{\theta} \sum_{i,j} e^{-\gamma d_{ij}}$$

Since  $e^{-\gamma d_{ij}}$  decays exponentially with  $d_{ij}$ , the sum is dominated by token pairs with small  $d_{ij}$ , corresponding to local interactions. Moreover, as the exponential decay ensures convergence of the sum, the total gradient norm remains bounded independently of sequence length. Consequently, parameter updates are primarily influenced by local token neighborhoods, contributions from distant tokens diminish exponentially, limiting their impact on parameter updates, and the bounded total gradient norm prevents gradient explosion, ensuring stable optimization dynamics.

In conclusion, our analysis of FourierRoFormer reveals its ability to approximate and interpret learned parameters. Our gradient analysis confirmed exponential decay with token distance, ensuring stable and localized training dynamics. These findings provide theoretical backing for the design of FourierRoFormer and its scalability to longer sequences.

## D ROPE COMPATIBILITY ANALYSIS

In this section we examine how the Fourier modulation in FourierRoFormer interacts with Rotary Position Embeddings (RoPE), and demonstrate that the combined attention mechanism retains key geometric properties of RoPE, including translation equivariance, relative position dependence, and structural decomposition.

**Theorem 3** (RoPE-Fourier Compatibility). *In FourierRoFormer, the modulated RoPE attention score* 

$$\mathbf{S}_{mn} = \frac{\langle \mathbf{R}_{\theta,m} \mathbf{q}_m, \mathbf{R}_{\theta,n} \mathbf{k}_n \rangle}{\sqrt{d}} \cdot \mathcal{M}(|m-n|) \cdot e^{-\gamma|m-n|}$$

is translation equivariant, depends only on relative positions, and admits a multiplicative decomposition. Specifically, for any shift  $\tau \in \mathbb{Z}$ , we have  $\mathbf{S}_{(m+\tau)(n+\tau)} = \mathbf{S}_{mn}$ , and  $\mathbf{S}_{mn}$  can be expressed as  $\mathbf{S}_{mn} = f(m-n, \mathbf{q}_m, \mathbf{k}_n)$  for some function f independent of absolute positions. Moreover, the score factorizes as  $\mathbf{S}_{mn} = \mathbf{S}_{mn}^{RoPE} \cdot \mathbf{S}_{mn}^{Fourier}$ , where  $\mathbf{S}_{mn}^{RoPE}$  is the standard RoPE attention score and  $\mathbf{S}_{mn}^{Fourier} = \mathcal{M}(|m-n|) \cdot e^{-\gamma|m-n|}$ .

*Proof.* We verify each property in turn. For translation equivariance, observe:

$$\mathbf{S}_{(m+\tau)(n+\tau)} = \frac{\langle \mathbf{R}_{\theta,m+\tau} \mathbf{q}_{m+\tau}, \mathbf{R}_{\theta,n+\tau} \mathbf{k}_{n+\tau} \rangle}{\sqrt{d}} \cdot \mathcal{M}(|m-n|) \cdot \mathcal{D}(|m-n|)$$

using  $|(m+\tau)-(n+\tau)|=|m-n|$ , and the RoPE invariance  $\mathbf{R}_{\theta,p+\tau}\mathbf{x}_{p+\tau}=\mathbf{R}_{\theta,p}\mathbf{x}_{p}$ . Hence,  $\mathbf{S}_{(m+\tau)(n+\tau)}=\mathbf{S}_{mn}$ . For relative position dependence, the RoPE inner product depends only on relative positions  $\langle \mathbf{R}_{\theta,m}\mathbf{q}_{m},\mathbf{R}_{\theta,n}\mathbf{k}_{n}\rangle=g(m-n,\mathbf{q}_{m},\mathbf{k}_{n})$  for some function g. Since  $\mathcal{M}$  and  $\mathcal{D}$  depend only on |m-n|, it follows that:

$$\mathbf{S}_{mn} = \frac{g(m-n, \mathbf{q}_m, \mathbf{k}_n)}{\sqrt{d}} \cdot \mathcal{M}(|m-n|) \cdot \mathcal{D}(|m-n|) = f(m-n, \mathbf{q}_m, \mathbf{k}_n)$$

For the decomposition, define:

$$\mathbf{S}_{mn}^{\text{RoPE}} = \frac{\langle \mathbf{R}_{\theta,m} \mathbf{q}_m, \mathbf{R}_{\theta,n} \mathbf{k}_n \rangle}{\sqrt{d}}, \quad \mathbf{S}_{mn}^{\text{Fourier}} = \mathcal{M}(|m-n|) \cdot \mathcal{D}(|m-n|)$$

Thus, by construction,  $\mathbf{S}_{mn} = \mathbf{S}_{mn}^{\text{RoPE}} \cdot \mathbf{S}_{mn}^{\text{Fourier}}$ .

To further understand the role of Fourier modulation, we observe that in the absence of learned Fourier components, FourierRoFormer simplifies to standard RoPE attention, as formalized below.

**Lemma 3** (RoPE Recovery). When all Fourier amplitudes  $a_k = 0$  or K = 0, FourierRoFormer reduces to standard RoPE attention with uniform modulation  $\mathcal{M}(d) = 0.5$ .

*Proof.* If  $a_k = 0$  for all k or equivalently K = 0, the modulation function simplifies to

$$\mathcal{M}(d) = \tanh(0) \cdot 0.5 + 0.5 = 0.5$$

Substituting into the attention score expression, we obtain

$$\mathbf{S}_{mn} = \frac{\langle \mathbf{R}_{\theta,m} \mathbf{q}_m, \mathbf{R}_{\theta,n} \mathbf{k}_n \rangle}{\sqrt{d}} \cdot 0.5 \cdot \mathcal{D}(|m-n|)$$

This corresponds to the standard RoPE attention, scaled by a constant factor and modulated by the damping function  $\mathcal{D}(|m-n|)$ . The structure of RoPE is thus preserved in the absence of active Fourier components.

Building on the compatibility and recovery properties established earlier, we conclude that FourierRoFormer extends RoPE by introducing learnable modulation while preserving its core structural advantages, as summarized in the following corollary.

**Corollary 3** (Enhanced Position Encoding). FourierRoFormer strictly enhances RoPE by preserving all of its beneficial properties, while introducing learnable frequency-based attention modulation and maintaining stable gradients through multiplicative interactions between the RoPE and Fourier components.

*Proof.* By Theorem 3, FourierRoFormer preserves the translation equivariance and relative position dependence of RoPE, ensuring that attention scores remain functions of relative positions only. Furthermore, the multiplicative decomposition of the attention score into a RoPE term and a Fourier modulation term preserves the structural properties of RoPE while introducing additional expressivity. Specifically, the Fourier modulation term  $\mathcal{M}(|m-n|)$  augments the standard RoPE attention with learnable, frequency-based modulation over token distances, enabling the model to adaptively emphasize or attenuate specific distance patterns. By Lemma 3, in the limiting case where  $a_k = 0$  for all k, FourierRoFormer recovers standard RoPE attention, confirming that RoPE is a special case within this generalized framework. Finally, the multiplicative interaction between the RoPE and Fourier terms maintains well-behaved gradients, as each component is bounded and differentiable, ensuring stable optimization. Therefore, FourierRoFormer strictly extends RoPE by preserving its key properties while enhancing its expressivity through learnable frequency modulation and maintaining stable training dynamics.

Building on Theorem 3, Lemma 3, and Corollary 3, FourierRoFormer generalizes RoPE by embedding its geometric properties within a learnable modulation framework. It preserves translation equivariance and relative position encoding, while enhancing expressivity through frequency-based modulation. This theoretical foundation highlights both the model's gradient stability and its adaptability to complex positional patterns.

#### E EXPERIMENTAL SETUP

All experiments are implemented in PYTORCH and executed on NVIDIA A40 GPUs with 48GB memory. To ensure fair comparison, we adopt a uniform training protocol, varying only key architectural hyperparameters. The *small*, *medium*, and *large* variants have embedding dimensions of 192, 384, and 576, respectively. The small and medium models use six attention heads, while the large model uses twelve. Transformer depth is six layers for the small model and twelve for the others.

Given the limited number of runs (n=5) and multiple comparisons across datasets, we adopt conservative statistical practices. We report confidence intervals alongside means and standard deviations. For significance testing, we use paired t-tests with Bonferroni correction across the 4 datasets tested, requiring p; 0.0125 for significance. We acknowledge that with 5 runs, detecting small effect sizes reliably is challenging, and focus our claims on improvements exceeding 2 percentage points.

**Baseline Methods and Comparisons:** We evaluate against three categories of methods: (1) Standard vision transformers (ViT, DeiT, RoFormer), (2) Recent positional encoding methods (ALiBi,

Context-aware Biases, Functional Interpolation), and (3) Spectral transformer methods (GFNet, WaveViT, SpectFormer, SVT).

**Relationship to Fourier Features.** Our approach differs fundamentally from coordinate-based Fourier features (Tancik et al., 2020), as detailed in table 10.

Table 10: Detailed comparison with Tancik et al. Fourier Features [26] highlighting fundamental differences in approach, application, and technical mechanism.

Aspect	Tancik et al. [26]	FourierRoFormer
<b>Application Domain</b>	Coordinate networks (NeRF, etc.)	Vision transformer attention
Target Problem	High-frequency function learning	Multi-scale attention modulation
Input Type	Continuous coordinates (x,y,z)	Discrete token sequences
Frequency Selection	Fixed random frequencies	Learnable adaptive frequencies
Parameter Learning	Static random $\gamma$ , fixed $\omega$	End-to-end learned $\{a_k, \omega_k, \phi_k\}$
Architecture Role	Input feature enhancement	Attention mechanism modulation
<b>Optimization Target</b>	Coordinate-to-value mapping	Token-to-token attention patterns
Data Dependency	Task-independent frequencies	Dataset-specific specialization
Interpretability	Fixed spectral bias	Learned frequency-pattern alignment
Scalability	Limited to coord. resolution	Scales with sequence length
<b>Evaluation Domain</b>	3D reconstruction, view synthesis	Image classification, detection
Core Innovation	Random Fourier input mapping	Learnable attention modulation

**Key Technical Distinctions:** Tancik et al. use fixed random frequencies for coordinate mapping, while we learn adaptive frequencies that specialize during training. Their method targets continuous coordinate functions, while ours operates on discrete token interactions. They enhance input representations, while we modulate attention mechanisms. Their approach uses static spectral bias, while ours learns dynamic patterns aligned with visual semantics.

Both methods leverage Fourier analysis but address fundamentally different problems: coordinate-based function approximation versus attention-based visual understanding.

**Spectral Transformer Baselines:** We include comprehensive comparisons with recent spectral methods: GFNet (Rao et al., 2021) uses fixed Fourier transforms for token mixing, while WaveViT (Yao et al., 2022) employs fixed wavelet transforms for multi-scale processing. Spect-Former (Patro et al., 2025b) provides a hybrid frequency-domain transformer with limited adaptability, and SVT (Oyallon et al., 2018) uses scattering-based spectral filtering with fixed wavelets.

**Key Differentiator:** Unlike these methods using fixed spectral transforms, FourierRoFormer learns adaptive frequency patterns  $\{a_k, \omega_k, \phi_k\}$  that specialize during training to capture dataset-specific visual patterns.

Memory requirements scaled with model complexity: small models required 11GB of GPU memory per run, medium models 18GB, and large models 32GB. Training times varied by dataset size and model scale: small models trained for approximately 5 hours on CIFAR-100, medium models for 12 hours, and large models for 22 hours. For ImageNet-subset, training times increased to 14, 28, and 48 hours respectively, while Oxford-Flowers102 required approximately 4, 9, and 17 hours for the three model sizes. The total compute for all experiments, including ablation studies and the 5 runs per configuration for statistical validation, amounted to approximately 2,100 GPU-hours. Inference overhead remains minimal, with the medium-sized FourierRoFormer processing 215 images/second on CIFAR-100 versus 220 for RoFormer on identical hardware. A detailed analysis of computational requirements for each dataset and model configuration is provided in Appendix E.1.

For CIFAR datasets, we use  $4\times 4$  image patches, while Oxford-Flowers102 and ImageNet use  $16\times 16$  patches. All models are trained with a batch size of 128 and optimized using AdamW with weight decay of 0.05. Learning rates follow a cosine decay schedule starting at  $5\times 10^{-4}$ , and models are trained for 20021 epochs. For ImageNet, standard data augmentation is used, including random resized crops and horizontal flips during training, and center cropping for evaluation.

Our DeiT implementation preserves the core architecture while adapting several components for fair comparison. We retain DeiT's training improvements such as strong regularization techniques but standardize the training duration to 200 epochs across all models rather than using the original

300+ epoch schedule. While maintaining the distillation token approach, we use a consistent teacher model across experiments. All optimization hyperparameters are aligned with our unified training protocol as described above, ensuring that performance differences arise primarily from architectural innovations rather than variations in training procedures.

Unless noted otherwise, FOURIERROFORMER is initialized with four learnable Fourier components, with frequencies linearly spaced between 0.1 and 2.0, an amplitude of 0.1, zero phase, and a damping coefficient of  $\gamma=0.01$ . This configuration ensures consistency across ablation studies, allowing performance differences to be directly attributed to the architectural choices under investigation.

#### E.1 COMPUTATIONAL RESOURCES

Our experimental framework was implemented in PyTorch and executed on NVIDIA A40 GPUs with 48GB of VRAM. Memory requirements scaled with model size: small models (192d, 6h, 6l) required 11GB memory with batch size 128, medium models (384d, 6h, 12l) used 18GB, and large models (576d, 12h, 12l) used 32GB. For the largest models on ImageNet-subset, we reduced the batch size to 64 to fit within memory constraints.

**Spectral Method Resource Comparison.** We conducted comprehensive resource analysis comparing FourierRoFormer with spectral transformer methods:

Table 11: Detailed resource comparison showing FourierRoFormer's superior resource efficiency compared to spectral transformer baselines.

Method	Memory	Peak Memory	Training Time	Energy (kWh)	CO <sub>2</sub> (kg)	Efficiency
RoFormer-M	18.0 GB	19.2 GB	12.0h	28.8	11.5	6.83
GFNet-H-B	21.5 GB	24.1 GB	16.8h	40.3	16.1	4.12
WaveViT-B	19.8 GB	22.4 GB	15.2h	36.5	14.6	5.46
SpectFormer-H-B	19.2 GB	21.8 GB	14.5h	34.8	13.9	5.89
SVT-H-B	19.5 GB	22.1 GB	15.8h	37.9	15.2	5.39
FourierRoFormer-M	18.1 GB	19.4 GB	12.3h	29.5	11.8	7.21
vs Best Spectral	-6.1%	-11.0%	-15.2%	-15.2%	-15.2%	+22.4%

**Resource Efficiency Metric:**  $\frac{\text{Top-1 Accuracy}^2}{\text{Training Time (h)} \times \text{Peak Memory (GB)}}$  captures accuracy-resource tradeoff.

#### F ABLATION STUDIES

We conduct comprehensive ablation studies to understand the contribution of each component in FourierRoFormer. All experiments in this section use the medium-sized model (384d, 6h, 12l) on CIFAR-100 unless otherwise specified.

**Quantitative Frequency Learning Validation.** We provide concrete empirical evidence that FourierRoFormer learns distinct frequency specialization during training. Table 12 shows quantitative tracking of frequency component evolution during ImageNet-1K training:

Table 12: Quantitative validation of frequency learning showing component specialization and correlation with visual patterns during ImageNet-1K training.

Component	Initial Amp	Final Amp	Learned Freq (Hz)	Visual Pattern	Correlation
k=1	$0.10 \pm 0.02$	0.43	0.3	Global object shape	r = 0.78
k=2	$0.10 \pm 0.02$	0.31	1.1	Object boundaries	r = 0.85
k=3	$0.10 \pm 0.02$	0.18	2.4	Fine textures	r = 0.71
k=4	$0.10 \pm 0.02$	0.08	3.2	Noise/artifacts	r = 0.34

**Three-Phase Training Dynamics.** Our analysis reveals distinct learning phases with measurable specialization metrics:

This quantitative analysis confirms that different frequency components learn to capture complementary visual patterns, with the strongest correlation (r = 0.85) achieved for object boundary detection at 1.1 Hz.

**Post-Attention vs Pre-Attention Modulation.** We provide comprehensive empirical validation for our design choice:

Table 13: Three-phase frequency learning progression with quantitative specialization metrics showing evolution from uniform exploration to structured hierarchy.

Phase	Epochs	Specialization $\sigma$	Coefficient Variation	Attention Entropy	Stability
Exploration	0-40	0.02	0.12	$3.41 \pm 0.18$	< 30%
Specialization	40-120	0.12	0.68	$3.38 \pm 0.12$	70%
Convergence	120+	0.31	0.91	$3.35 \pm 0.08$	>95%

Table 14: Comprehensive comparison of modulation placement showing superior performance and stability of post-attention design.

Modulation	ImageNet CIFAR Top-1 -100		Gradient $\sigma$	Convergence	Semantic Preservation	Training Stability	
Pre-attention Post-attention	82.3% <b>84.1%</b>	82.8% <b>84.26</b> %	0.41 <b>0.12</b>	Epoch 145 <b>Epoch 128</b>	0.72 <b>0.89</b>	Unstable <b>Stable</b>	
Improvement	+1.8pp	+1.46pp	-71%	-12%	+24%	Qualitative	

**Multi-Head Frequency Specialization.** When allowing head-specific frequency parameters, we observe emergent specialization:

Table 15: Multi-head frequency specialization showing automatic division of labor across attention heads with quantitative metrics.

Configuration	ImageNet Top-1	Head Group	Freq Range (Hz)	Attention Range	Energy % Energy %	Specialization Timeline
Uniform	84.1%	All heads	0.5-1.5	45 tokens	100%	None
Head-specific	84.6%	Heads 1-2 Heads 3-4 Heads 5-6	0.2-0.6 0.6-1.4 1.4-3.2	89 tokens 43 tokens 21 tokens	35% 40% 25%	Epoch 35 Epoch 42 Epoch 38

**Fourier Components and Damping.** We analyze the impact of each component by selective ablation, as shown in Table 16.

Fourier Components and Damping. We analyze the impact of each component by selective ablation, as shown in Table 16. Fourier modulation alone provides improvement (+4.43pp) over the RoFormer baseline, while damping alone contributes +2.09pp. When combined, these components achieve a complementary effect, yielding +5.84pp total improvement. Our experiments with varying the number of Fourier components (K) show that 4-8 components provides the optimal balance between expressivity and overfitting, with K=8 achieving the best performance (+6.53pp). Similarly, moderate damping ( $\gamma$ =0.01) yields the best results among the damping coefficients tested.

**Frequency Initialization Strategies.** We also investigate different approaches for initializing the Fourier component frequencies, as shown in Table 17. Logarithmic spacing achieves the best performance (84.62%), providing better coverage across the frequency spectrum compared to linear spacing. Random initialization performs worse (83.91%), suggesting that a structured approach to frequency initialization aids optimization. Low-frequency bias initialization shows moderate performance, indicating that while low frequencies are important, a balanced coverage across the spectrum is more effective.

#### G COMPUTATIONAL COMPLEXITY ANALYSIS

For completeness, we analyze the computational overhead introduced by the Fourier modulation components in FourierRoFormer. Let n denote the input sequence length, d the feature dimension, and  $\kappa$  the number of Fourier components. The computation of the Fourier modulation function requires evaluating  $\kappa$  cosine terms for each token pair, computing the modulation, and applying non-linear scaling. Since there are  $\mathcal{O}(n^2)$  token pairs in the attention mechanism Vaswani et al. (2017), this results in an overall computational cost of  $\mathcal{O}(\kappa n^2)$  operations Rahimi & Recht (2008); Tancik et al. (2020).

**Comprehensive Efficiency Comparison with Spectral Methods.** We provide detailed efficiency analysis comparing FourierRoFormer with spectral transformer baselines:

Table 16: Comprehensive ablation study on CIFAR-100 showing complementary benefits of components.

Configuration	Accuracy (%)	$\Delta$ vs RoFormer	Params (M)	GFLOPs			
RoFormer (baseline)	78.42	-	24.75	4.60			
+ Fourier only	82.85	+4.43	24.75	4.61			
+ Damping only	80.51	+2.09	24.75	4.60			
+ Both (Full model)	84.26	+5.84	24.76	4.63			
Fourier Component Variations							
K=2 components	82.54	+4.12	24.75	4.61			
K=4 components 84.26		+5.84	24.76	4.63			
K=8 components 84.95		+6.53	24.76	4.63			
K=16 components	84.72	+6.30	24.77	4.64			
Damping Coefficient Analysis							
$\gamma = 0.001$	83.45	+5.03 24.76		4.63			
$\gamma = 0.01$	84.26	+5.84	24.76	4.63			
$\gamma = 0.05$	83.87	+5.45	24.76	4.63			
$\dot{\gamma} = 0.1$	82.93	+4.51	24.76	4.63			

Table 17: Comparison of frequency initialization strategies on CIFAR-100.

Accuracy (%)	Description
84.26	Frequencies evenly spaced 0.1-2.0
84.62	Log-spaced frequencies
83.91	Random frequencies 0.1-2.0
84.08	Emphasis on low frequencies
	84.26 <b>84.62</b> 83.91

## **Efficiency Metrics Defined:**

- Efficiency Score =  $\frac{\text{Top-1 Accuracy}}{\log(\text{Params}) \times \sqrt{\text{Training Time}}}$  (higher is better)
- Parameter Efficiency =  $\frac{\text{Top-1 Accuracy}}{\text{Params (M)}}$  (accuracy per million parameters)
- Computational Efficiency =  $\frac{\text{Top-1 Accuracy}}{\text{FLOPs (G)}}$  (accuracy per GFLOP)

**Key Findings:** The approach introduces minimal overhead with only 0.04% parameter increase and 0.7% FLOPs increase over RoFormer. It achieves a superior tradeoff with 23% better efficiency score than the best spectral baseline while using 25% fewer parameters. The method provides practical advantage by maintaining standard transformer architecture compatibility unlike spectral methods requiring architectural overhaul.

The additional computational cost of FourierRoFormer compared to standard ViT or RoFormer is minimal, with only 0.01M additional parameters (0.04%) from the learnable Fourier components. During inference, FourierRoFormer processes approximately 215 images/second on our medium model configuration for CIFAR-100, compared to 220 images/second for RoFormer and 218 images/second for standard ViT on identical hardware, demonstrating negligible runtime overhead for improved accuracy gains.

Table 18: Comprehensive efficiency analysis showing FourierRoFormer achieves optimal accuracy-efficiency tradeoff compared to spectral transformer methods.

Method	Params (M)	Memory (GB)	Throughput (img/s)	Training Time (h)	FLOPs (G)	Top-1 (%)	Efficiency Score	Parameter Efficiency	
RoFormer-M	24.75	18.0	220	12.0	4.60	81.9	3.33	3.31	
GFNet-H-B	54.0	21.5	185	16.8	8.6	82.9	2.41	1.54	
WaveViT-B	33.5	19.8	195	15.2	6.8	84.8	2.98	2.53	
SpectFormer-H-B	33.1	19.2	195	14.5	6.3	85.1	3.21	2.57	
SVT-H-B	32.8	19.5	190	15.8	6.5	85.2	3.18	2.60	
FourierRoFormer-M	24.76	18.1	215	12.3	4.63	84.1	3.91	3.40	
Efficiency Advantage vs Best Spectral Baseline (SVT-H-B)									
Relative Advantage	-24.5%	-7.2%	+13.2%	-22.2%	-28.8%	-1.1pp	+23%	+31%	